

On the Chain Pair Simplification Problem

Chenglin Fan ^{*} Omrit Filtser [†] Matthew J. Katz [‡] Tim Wylie [§]
 Binhai Zhu [¶]

Abstract

The problem of efficiently computing and visualizing the structural resemblance between a pair of protein backbones in 3D has led Bereg et al. [BJW⁺08] to pose the Chain Pair Simplification problem (CPS). In this problem, given two polygonal chains A and B of lengths m and n , respectively, one needs to simplify them simultaneously, such that each of the resulting simplified chains, A' and B' , is of length at most k and the discrete Fréchet distance between A' and B' is at most δ , where k and δ are given parameters.

In this paper we study the complexity of CPS under the discrete Fréchet distance (CPS-3F), i.e., where the quality of the simplifications is also measured by the discrete Fréchet distance. Since CPS-3F was posed in 2008, its complexity has remained open. However, it was believed to be **NP**-complete, since CPS under the Hausdorff distance (CPS-2H) was shown to be **NP**-complete. We first prove that the weighted version of CPS-3F is indeed weakly **NP**-complete even on the line, based on a reduction from the set partition problem. Then, we prove that CPS-3F is actually polynomially solvable, by presenting an $O(m^2n^2 \min\{m, n\})$ time algorithm for the corresponding minimization problem. In fact, we prove a stronger statement, implying, for example, that if weights are assigned to the vertices of only one of the chains, then the problem remains polynomially solvable. We also study a few less rigid variants of CPS and present efficient solutions for them.

Finally, we present some experimental results that suggest that (the minimization version of) CPS-3F is significantly better than previous algorithms for the motivating biological application.

1 Introduction

Polygonal curves play an important role in many applied areas, such as 3D modeling in computer vision, map matching in GIS, and protein backbone structural alignment and

^{*}Montana State University, Bozeman, MT, 59717-3880 USA; chenglin.fan@msu.montana.edu

[†]Ben-Gurion University of the Negev, Beer-Sheva 84105, Israel; omritna@post.bgu.ac.il

[‡]Ben-Gurion University of the Negev, Beer-Sheva 84105, Israel; matya@cs.bgu.ac.il

[§]The University of Texas-Pan American, Edinburg, TX, 78539 USA; wylie@utpa.edu

[¶]Montana State University, Bozeman, MT, 59717-3880 USA; bhz@cs.montana.edu

comparison in computational biology. Many different methods exist to compare curves in these (and in many other) applications, where one of the more prevalent methods is the Fréchet distance [Fré06].

The *Fréchet distance* is often described by an analogy of a man and a dog connected by a leash, each walking along a curve from its starting point to its end point. Both the man and the dog can control their speed but they are not allowed to backtrack. The Fréchet distance between the two curves is the minimum length of a leash that is sufficient for traversing both curves in this manner.

The *discrete Fréchet distance* is a simpler version, where, instead of continuous curves, we are given finite sequences of points, obtained, e.g., by sampling the continuous curves, or corresponding to the vertices of polygonal chains. Now, the man and the dog only hop monotonically along the sequences of points. The discrete Fréchet distance is considered a good approximation of the continuous distance.

One promising application of the discrete Fréchet distance has been protein backbone comparison. Within structural biology, polygonal curve alignment and comparison is a central problem in relation to proteins. Proteins are usually studied using RMSD (Root Mean Square Deviation), but recently the discrete Fréchet distance was used to align and compare protein backbones, which yielded favorable results in many instances [JXZ08, WLZ11]. In this application, the discrete version of the Fréchet distance makes more sense, because by using it the alignment is done with respect to the vertices of the chains, which represent α -carbon atoms. Applying the continuous Fréchet distance will result in mapping of arbitrary points, which is not meaningful biologically.

There may be as many as 500~600 α -carbon atoms along a protein backbone, which are the nodes (i.e., points) of our chain. This makes efficient computation essential, and is one of the reasons for considering simplification. In general, given a chain A of n vertices, a simplification of A is a chain A' such that A' is “close” to A and the number of vertices in A' is significantly less than n . The problem of simplifying a 3D polygonal chains under the discrete Fréchet distance was first addressed by Bereg et al. [BJW⁺08].

Simplifying two aligned chains independently does not necessarily preserve the resemblance between the chains; see Figure 1. Thus, the following question arises: Is it possible to simplify both chains in a way that will retain the resemblance between them? This question has led Bereg et al. [BJW⁺08] to pose the Chain Pair Simplification problem (CPS). In this problem, the goal is to simplify both chains simultaneously, so that the discrete Fréchet distance between the resulting simplifications is bounded. More precisely, given two chains A and B of lengths m and n , respectively, an integer k and three real numbers $\delta_1, \delta_2, \delta_3$, one needs to find two chains A', B' with vertices from A, B , respectively, each of length at most k , such that $d_1(A, A') \leq \delta_1$, $d_2(B, B') \leq \delta_2$, $d_{dF}(A', B') \leq \delta_3$ (d_1 and d_2 can be any similarity measures and d_{dF} is the discrete Fréchet distance). When the chains are simplified using the Hausdorff distance, i.e., d_1, d_2 is the Hausdorff distance (CPS-2H), the problem becomes **NP**-complete [BJW⁺08]. However, the complexity of the version in which d_1, d_2 is the discrete Fréchet distance (CPS-3F) has been open since 2008.

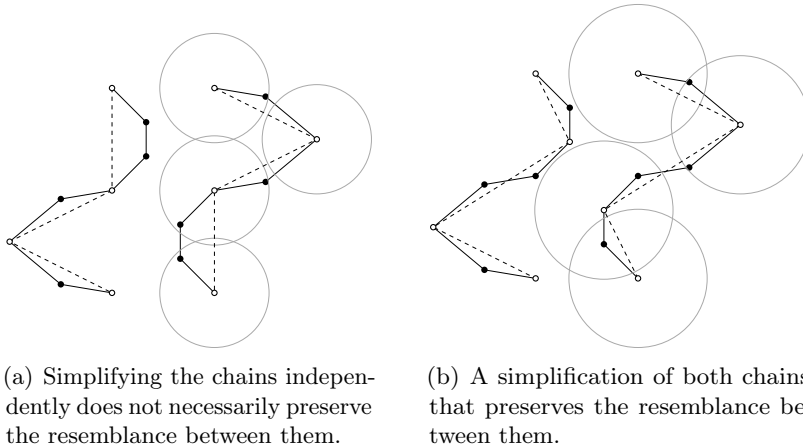


Figure 1: Independent simplification vs. simultaneous simplification. Each chain simplification consists of 4 vertices (marked by empty circles) chosen from the corresponding chain. The unit disks illustrate the Fréchet distance between the right chain in each of the figures and its corresponding simplification; their radius in (b) is larger.

Related work. The Fréchet distance and its variants have been studied extensively in the past two decades. Alt and Godau [AG95] gave an $O(mn \log mn)$ -time algorithm for computing the Fréchet distance between two polygonal curves of lengths m and n . This result in the plane was recently improved by Buchin et al [BBMM14]. The discrete Fréchet distance was originally defined by Eiter and Mannila [EM94], who also presented an $O(mn)$ -time algorithm for computing it. A slightly sub-quadratic algorithm was given recently by Agarwal et al. [AAKS14].

As mentioned earlier, Bereg et al. [BJW⁺08] were the first to study simplification problems under the discrete Fréchet distance. They considered two such problems. In the first, the goal is to minimize the number of vertices in the simplification, given a bound on the distance between the original chain and its simplification, and, in the second problem, the goal is to minimize this distance, given a bound k on the number of vertices in the simplification. They presented an $O(n^2)$ -time algorithm for the former problem and an $O(n^3)$ -time algorithm for the latter problem, both using dynamic programming, for the case where the vertices of the simplification are from the original chain. (For the arbitrary vertices case, they solve the problems in $O(n \log n)$ time and in $O(kn \log n \log(n/k))$ time, respectively.) Driemel and Har-Peled [DH13] showed how to preprocess a polygonal curve in near-linear time and space, such that, given an integer $k > 0$, one can compute a simplification in $O(k)$ time which has $2k - 1$ vertices of the original curve and is optimal up to a constant factor (w.r.t. the continuous Fréchet distance), compared to any curve consisting of k arbitrary vertices.

For the chain pair simplification problem (CPS), Bereg et al. [BJW⁺08] proved that CPS-2H is **NP**-complete, and conjectured that so is CPS-3F. Wylie et al. [WLZ11] gave a heuristic algorithm for CPS-3F, using a greedy method with backtracking, and based on the assumption that the (Euclidean) distance between adjacent α -carbon atoms in a protein backbone is almost fixed. More recently, Wylie and Zhu [WZ13] presented an approximation algorithm with approximation ratio 2 for the optimization version of CPS-3F. Their algorithm actually solves the optimization version of a related problem called CPS-3F⁺, it uses dynamic programming and its running time is between $O(mn)$ and $O(m^2n^2)$ depending on the input simplification parameters.

Some special cases of CPS-3F have recently been studied. Motivated by the need to reduce sensitivity to outliers when comparing curves, Ben Avraham et al. [AFK⁺14] studied the discrete Fréchet distance with shortcuts problem. In the one-sided variant, the dog is allowed to jump to any point that comes later in its sequence, rather than just to the next point. The man has to visit the points in its sequence, one after the other, as in the standard discrete Fréchet distance problem. In the two-sided variant, both the man and the dog are allowed to skip points. Unlike CPS-3F, the difference between an original chain and its simplification (in the two-sided variant) can be big, since the sole goal is to minimize the discrete Fréchet distance between the two simplified chains. (For this reason, Ben Avraham et al. do not allow both the man and the dog to move simultaneously, since, otherwise, they would both jump directly to their final points.) Moreover, the length of a simplification is only bounded by the length of the corresponding chain. Both variants of the shortcuts problem can be solved in subquadratic time.

The one-sided variant of the (continuous) Fréchet distance with shortcuts problem was studied by Driemel et al. [DH13], who considered the problem assuming the curves are *c-packed* and shortcuts start and end at vertices of the noisy curve. They gave a near-linear time $(3 + \varepsilon)$ -approximation algorithm. Buchin et al. [BDS14] proved that the more general variant, where shortcuts can be taken at any point along the noisy curve, is **NP**-hard, and gave an $O(n^3 \log n)$ -time 3-approximation algorithm for the corresponding decision problem. Another approach for handling outliers (that is still somewhat related to our work) was proposed by Buchin et al. [BBW09], who studied the partial curve matching problem under the (continuous) Fréchet distance. That is, given two curves and a threshold δ , find subcurves of maximum total length that are close to each other w.r.t. δ .

Our results. In Section 3 we introduce the weighted chain pair simplification problem and prove that weighted CPS-3F is weakly **NP**-complete. In Section 4, we resolve the question concerning the complexity of CPS-3F by proving that it is polynomially solvable, contrary to what was believed. We do this by presenting a polynomial-time algorithm for the corresponding optimization problem. We actually prove a stronger statement, implying, for example, that if weights are assigned to the vertices of only one of the chains, then the problem remains polynomially solvable. In Section 5 we devise a sophisticated $O(m^2n^2 \min\{m, n\})$ -time dynamic programming algorithm for the minimization problem of

CPS-3F. Besides being interesting from a theoretical point of view, only after developing (and implementing) this algorithm, were we able to apply the CPS-3F minimization problem to datasets from the Protein Data Bank (PDB), see below.

In section 6 we study several less rigid variants of CPS-3F. In particular, we improve the result of Bereg et al. [BJW⁺08] mentioned above on the problem of finding the best simplification of a given length under the discrete Fréchet distance, by presenting a more general $O(n^2 \log n)$ -time algorithm (rather than an $O(n^3)$ -time algorithm).

Finally, in Section 7 we present some empirical results comparing (the minimization problems of) CPS-3F⁺ [WZ13] (the best available algorithm prior to this work) and CPS-3F using datasets from the PDB, and showing that with the latter we get much smaller simplifications (obeying the same distance bounds).

2 Preliminaries

Let $A = (a_1 \dots, a_m)$ and $B = (b_1, \dots, b_n)$ be two sequences of m and n points, respectively, in \mathbb{R}^k . The discrete Fréchet distance $d_{dF}(A, B)$ between A and B is defined as follows. Fix a distance $\delta > 0$ and consider the Cartesian product $A \times B$ as the vertex set of a directed graph G_δ whose edge set is

$$\begin{aligned} E_\delta = & \{((a_i, b_j), (a_{i+1}, b_j)) \mid d(a_i, b_j), d(a_{i+1}, b_j) \leq \delta\} \cup \\ & \{((a_i, b_j), (a_i, b_{j+1})) \mid d(a_i, b_j), d(a_i, b_{j+1}) \leq \delta\} \cup \\ & \{((a_i, b_j), (a_{i+1}, b_{j+1})) \mid d(a_i, b_j), d(a_{i+1}, b_{j+1}) \leq \delta\}. \end{aligned}$$

Then $d_{dF}(A, B)$ is the smallest $\delta > 0$ for which (a_m, b_n) is reachable from (a_1, b_1) in G_δ .

The chain pair simplification problem (CPS) is formally defined as follows.

Problem 1 (Chain Pair Simplification).

Instance: Given a pair of polygonal chains A and B of lengths m and n , respectively, an integer k , and three real numbers $\delta_1, \delta_2, \delta_3 > 0$.

Problem: Does there exist a pair of chains A', B' each of at most k vertices, such that the vertices of A', B' are from A, B , respectively, and $d_1(A, A') \leq \delta_1$, $d_2(B, B') \leq \delta_2$, and $d_{dF}(A', B') \leq \delta_3$?

When $d_1 = d_2 = d_H$, the problem is **NP**-complete and is called CPS-2H, and when $d_1 = d_2 = d_{dF}$, the problem is called CPS-3F.

3 Weighted Chain Pair Simplification (WCPS-3F)

We first introduce and consider a more general version of CPS-3F, namely, Weighted CPS-3F. In the weighted version of the chain pair simplification problem, the vertices of the

chains A and B are assigned arbitrary weights, and, instead of limiting the length of the simplifications, one limits their weights. That is, the total weight of each simplification must not exceed a given value. The problem is formally defined as follows.

Problem 2 (Weighted Chain Pair Simplification).

Instance: Given a pair of 3D chains A and B , with lengths m and n , respectively, an integer k , three real numbers $\delta_1, \delta_2, \delta_3 > 0$, and a weight function $C : \{a_1, \dots, a_m, b_1, \dots, b_n\} \rightarrow \mathbb{R}^+$.

Problem: Does there exist a pair of chains A', B' with $C(A'), C(B') \leq k$, such that the vertices of A', B' are from A, B respectively, $d_1(A, A') \leq \delta_1$, $d_2(B, B') \leq \delta_2$, and $d_{dF}(A', B') \leq \delta_3$?

When $d_1 = d_2 = d_{dF}$, the problem is called WCPS-3F. When $d_1 = d_2 = d_H$, the problem is **NP**-complete, since the non-weighted version (i.e., CPS-2H) is already **NP**-complete [BJW⁺08].

We prove that WCPS-3F is weakly **NP**-complete via a reduction from the *set partition* problem: Given a set of positive integers $S = \{s_1, \dots, s_n\}$, find two sets $P_1, P_2 \subset S$ such that $P_1 \cap P_2 = \emptyset$, $P_1 \cup P_2 = S$, and the sum of the numbers in P_1 equals the sum of the numbers in P_2 . This is a weakly **NP**-complete special case of the classic subset-sum problem.

Our reduction builds two curves with weights reflecting the values in S . We think of the two curves as the subsets of the partition of S . Although our problem requires positive weights, we also allow zero weights in our reduction for clarity. Later, we show how to remove these weights by slightly modifying the construction.

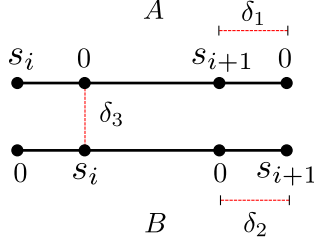


Figure 2: The reduction for the weighted chain pair simplification problem under the discrete Fréchet distance.

Theorem 1. *The weighted chain pair simplification problem under the discrete Fréchet distance is weakly **NP**-complete.*

Proof. Given the set of positive integers $S = \{s_1, \dots, s_n\}$, we construct two curves A and B in the plane, each of length $2n$. We denote the weight of a vertex x_i by $w(x_i)$. A is constructed as follows. The i 'th odd vertex of A has weight s_i , i.e. $w(a_{2i-1}) = s_i$, and coordinates $a_{2i-1} = (i, 1)$. The i 'th even vertex of A has coordinates $a_{2i} = (i + 0.2, 1)$ and weight zero. Similarly, the i 'th odd vertex of B has weight zero and coordinates $b_{2i-1} = (i, 0)$,

and the i 'th even vertex of B has coordinates $b_{2i} = (i + 0.2, 0)$ and weight s_i , i.e. $w(b_{2i}) = s_i$. Figure 2 depicts the vertices $a_{2i-1}, a_{2i}, a_{2(i+1)-1}, a_{2(i+1)}$ of A and $b_{2i-1}, b_{2i}, b_{2(i+1)-1}, b_{2(i+1)}$ of B . Finally, we set $\delta_1 = \delta_2 = 0.2$, $\delta_3 = 1$, and $k = \mathfrak{S}$, where \mathfrak{S} denotes the sum of the elements of S (i.e., $\mathfrak{S} = \sum_{j=1}^n s_j$).

We claim that S can be partitioned into two subsets, each of sum $\mathfrak{S}/2$, if and only if A and B can be simplified with the constraints $\delta_1 = \delta_2 = 0.2$, $\delta_3 = 1$ and $k = \mathfrak{S}/2$, i.e., $C(A'), C(B') \leq \mathfrak{S}/2$.

First, assume that S can be partitioned into sets S_A and S_B , such that $\sum_{s \in S_A} s = \sum_{s \in S_B} s = \mathfrak{S}/2$. We construct simplifications of A and of B as follows.

$$A' = \{a_{2i-1} \mid s_i \in S_A\} \cup \{a_{2i} \mid s_i \notin S_A\} \text{ and } B' = \{b_{2i} \mid s_i \in S_B\} \cup \{b_{2i-1} \mid s_i \notin S_B\}.$$

It is easy to see that $C(A'), C(B') \leq \mathfrak{S}/2$. Also, since $\{S_A, S_B\}$ is a partition of S , exactly one of the following holds, for any $1 \leq i \leq n$:

1. $a_{2i-1} \in A', b_{2i-1} \in B'$ and $a_{2i} \notin A', b_{2i} \notin B'$.
2. $a_{2i-1} \notin A', b_{2i-1} \notin B'$ and $a_{2i} \in A', b_{2i} \in B'$.

This implies that $d_{dF}(A, A') \leq 0.2 = \delta_1$, $d_{dF}(B, B') \leq 0.2 = \delta_2$ and $d_{dF}(A', B') \leq 1 = \delta_3$.

Now, assume there exist simplifications A', B' of A, B , such that $d_{dF}(A, A') \leq \delta_1 = 0.2$, $d_{dF}(B, B') \leq \delta_2 = 0.2$, $d_{dF}(A', B') \leq \delta_3 = 1$, and $C(A'), C(B') \leq k = \mathfrak{S}/2$. Since $\delta_1 = \delta_2 = 0.2$, for any $1 \leq i \leq n$, the simplification A' must contain one of a_{2i-1}, a_{2i} , and the simplification B' must contain one of b_{2i-1}, b_{2i} . Since $\delta_3 = 1$, for any i , at least one of the following two conditions holds: $a_{2i-1} \in A'$ and $b_{2i-1} \in B'$ or $a_{2i} \in A'$ and $b_{2i} \in B'$. Therefore, for any i , either $a_{2i-1} \in A$ or $b_{2i} \in B$, implying that s_i participates in either $C(A')$ or $C(B')$. However, since $C(A'), C(B') \leq \mathfrak{S}/2$, s_i cannot participate in both $C(A')$ and $C(B')$. It follows that $C(A') = C(B') = \mathfrak{S}/2$, and we get a partition of S into two sets, each of sum $\mathfrak{S}/2$.

Finally, we note that WCPS-3F is in **NP**. For an instance I with chains A, B , given simplifications A', B' , we can verify in polynomial time that $d_{dF}(A, A') \leq \delta_1$, $d_{dF}(B, B') \leq \delta_2$, $d_{dF}(A', B') \leq \delta_3$, and $C(A'), C(B') \leq k$. \square

Although our construction of A' and B' uses zero weights, a simple modification enables us to prove that the problem is weakly **NP**-complete also when only positive integral weights are allowed. Increase all the weights by 1, that is, $w(a_{2i-1}) = w(b_{2i}) = s_i + 1$ and $w(a_{2i}) = w(b_{2i-1}) = 1$, for $1 \leq i \leq n$, and set $k = \mathfrak{S}/2 + n$. It is easy to verify that our reduction still works. Finally, notice that we could overlay the two curves choosing $\delta_3 = 0$ and prove that the problem is still weakly **NP**-complete in one dimension.

4 Chain Pair Simplification (CPS-3F)

We now turn our attention to CPS-3F, which is the special case of WCPS-3F where each vertex has weight one.

We present an algorithm for the minimization version of CPS-3F. That is, we compute the minimum integer k^* , such that there exists a “walk”, as above, in which each of the dogs makes at most k^* hops. The answer to the decision problem is “yes” if and only if $k^* < k$.

Returning to the analogy of the man and the dog, we can extend it as follows. Consider a man and his dog connected by a leash of length δ_1 , and a woman and her dog connected by a leash of length δ_2 . The two dogs are also connected to each other by a leash of length δ_3 . The man and his dog are walking on the points of a chain A and the woman and her dog are walking on the points of a chain B . The dogs may skip points. The problem is to determine whether there exists a “walk” of the man and his dog on A and the woman and her dog on B , such that each of the dogs steps on at most k points.

Overview of the algorithm. We say that (a_i, a_p, b_j, b_q) is a *possible* configuration of the man, woman and the two dogs on the paths A and B , if $d(a_i, a_p) \leq \delta_1$, $d(b_j, b_q) \leq \delta_2$ and $d(a_p, b_q) \leq \delta_3$. Notice that there are at most $m^2 n^2$ such configurations. Now, let G be the DAG whose vertices are the possible configurations, such that there exists a (directed) edge from vertex $u = (a_i, a_p, b_j, b_q)$ to vertex $v = (a_{i'}, a_{p'}, b_{j'}, b_{q'})$ if and only if our gang can move from configuration u to configuration v . That is, if and only if $i \leq i' \leq i + 1$, $p \leq p'$, $j \leq j' \leq j + 1$, and $q \leq q'$. Notice that there are no cycles in G because backtracking is forbidden. For simplicity, we assume that the first and last points of A' (resp., of B') are a_1 and a_m (resp., b_1 and b_n), so the initial and final configurations are $s = (a_1, a_1, b_1, b_1)$ and $t = (a_m, a_m, b_n, b_n)$, respectively. (It is easy, however, to adapt the algorithm below to the case where the initial and final points of A' and B' are not specified, see remark below.) Our goal is to find a path from s to t in G . However, we want each of our dogs to step on at most k points, so, instead of searching for any path from s to t , we search for a path that minimizes the value $\max\{|A'|, |B'|\}$, and then check if this value is at most k .

For each edge $e = (u, v)$, we assign two weights, $w_A(e), w_B(e) \in \{0, 1\}$, in order to compute the number of hops in A' and in B' , respectively. $w_A(u, v) = 1$ if and only if the first dog jumps to a new point between configurations u and v (i.e., $p < p'$), and, similarly, $w_B(u, v) = 1$ if and only if the second dog jumps to a new point between u and v (i.e., $q < q'$). Thus, our goal is to find a path P from s to t in G , such that $\max\{\sum_{e \in P} w_A(e), \sum_{e \in P} w_B(e)\}$ is minimized.

Assume w.l.o.g. that $m \leq n$. Since $|A'| \leq m$ and $|B'| \leq n$, we maintain, for each vertex v of G , an array $X(v)$ of size m , where $X(v)[r]$ is the minimum number z such that v can be reached from s with (at most) r hops of the first dog and z hops of the second dog. We can construct these arrays by processing the vertices of G in topological order (i.e., a vertex is processed only after all its predecessors have been processed). This yields an algorithm of running time $O(m^3 n^3 \min\{m, n\})$, as described in Algorithm 1.

Running time. The number of vertices in G is $|V| = O(m^2 n^2)$. By the construction of the graph, for any vertex (a_i, a_p, b_j, b_q) the maximum number of outgoing edges is $O(mn)$.

Algorithm 1 CPS-3F

1. Create a directed graph $G = (V, E)$ with two weight functions w_A, w_B , such that:
 - V is the set of all configurations (a_i, a_p, b_j, b_q) with $d(a_i, a_p) \leq \delta_1$, $d(b_j, b_q) \leq \delta_2$, and $d(a_p, b_q) \leq \delta_3$.
 - $E = \{((a_i, a_p, b_j, b_q), (a_{i'}, a_{p'}, b_{j'}, b_{q'})) \mid i \leq i' \leq i+1, p \leq p', j \leq j' \leq j+1, q \leq q'\}$.
 - For each $((a_i, a_p, b_j, b_q), (a_{i'}, a_{p'}, b_{j'}, b_{q'})) \in E$, set
 - $w_A((a_i, a_p, b_j, b_q), (a_{i'}, a_{p'}, b_{j'}, b_{q'})) = \begin{cases} 1, & p < p' \\ 0, & \text{otherwise} \end{cases}$
 - $w_B((a_i, a_p, b_j, b_q), (a_{i'}, a_{p'}, b_{j'}, b_{q'})) = \begin{cases} 1, & q < q' \\ 0, & \text{otherwise} \end{cases}$
 2. Sort V topologically.
 3. Initialize the array $X(s)$ (i.e., set $X(s)[r] = 0$, for $r = 0, \dots, m-1$).
 4. For each $v \in V \setminus \{s\}$ (advancing from left to right in the sorted sequence) do:
 - (a) Initialize the array $X(v)$ (i.e., set $X(v)[r] = \infty$, for $r = 0, \dots, m-1$).
 - (b) For each r between 0 and $m-1$, compute $X(v)[r]$:
$$X(v)[r] = \min_{(u,v) \in E} \begin{cases} X(u)[r] + w_B(u, v), & w_A(u, v) = 0 \\ X(u)[r-1] + w_B(u, v), & w_A(u, v) = 1 \end{cases}$$
 5. Return $k^* = \min_r \max\{r, X(t)[r]\}$.
-

So we have $|E| = O(|V|mn) = O(m^3n^3)$. Thus, constructing the graph G in Step 1 takes $O(n^3m^3)$ time. Step 2 takes $O(|E|)$ time, while Step 3 takes $O(m)$ time. In Step 4, for each vertex v and for each index r , we consider all configurations that can directly precede v . So each edge of G participates in exactly m minimum computations, implying that Step 4 takes $O(|E|m)$ time. Step 5 takes $O(m)$ time. Thus, the total running time of the algorithm is $O(m^4n^3)$.

Theorem 2. *The chain pair simplification problem under the discrete Fréchet distance (CPS-3F) is polynomial, i.e., CPS-3F $\in \mathbf{P}$.*

Remark 1. As mentioned, we have assumed that the first and last points of A' (resp., B') are a_1 and a_m (resp., b_1 and b_n), so we have a single initial configuration (i.e., $s = (a_1, a_1, b_1, b_1)$) and a single final configuration (i.e., $t = (a_m, a_m, b_n, b_n)$). However, it is easy to adapt our algorithm to the case where the first and last points of the chains A' and B' are not specified. In this case, any possible configuration of the form (a_1, a_p, b_1, b_q) is considered a potential initial configuration, and any possible configuration of the form (a_m, a_p, b_n, b_q) is

considered a potential final configuration, where $1 \leq p \leq m$ and $1 \leq q \leq n$. Let S and T be the sets of potential initial and final configurations, respectively. (Then, $|S| = O(mn)$ and $|T| = O(mn)$.) We thus remove from G all edges entering a potential initial configuration, so that each such configuration becomes a “root” in the (topologically) sorted sequence. Now, in Step 3 we initialize the arrays of each $s \in S$ in total time $O(m^2n)$, and in Step 4 we only process the vertices that are not in S . The value $X(v)[r]$ for such a vertex v is now the minimum number z such that v can be reached from s with r hops of the first dog and z hops of the second dog, over *all* potential initial configurations $s \in S$. In the final step of the algorithm, we calculate the value k^* in $O(m)$ time, for each potential final configuration $t \in T$. The smallest value obtained is then the desired value. Since the number of potential final configurations is only $O(mn)$, the total running time of the final step of the algorithm is only $O(m^2n)$, and the running time of the entire algorithm remains $O(m^4n^3)$.

4.1 The weighted version

Weighted CPS-3F, which was shown to be weakly **NP**-complete in the previous section, can be solved in a similar manner, albeit with running time that depends on the number of different point weights in chain A (alternatively, B). We now explain how to adapt our algorithm to the weighted case. We first redefine the weight functions w_A and w_B (where $C(x)$ is the weight of point x):

$$\begin{aligned} \bullet \quad w_A((a_i, a_p, b_j, b_q), (a_{i'}, a_{p'}, b_{j'}, b_{q'})) &= \begin{cases} C(a_{p'}), & p < p' \\ 0, & \text{otherwise} \end{cases} \\ \bullet \quad w_B((a_i, a_p, b_j, b_q), (a_{i'}, a_{p'}, b_{j'}, b_{q'})) &= \begin{cases} C(b_{q'}), & q < q' \\ 0, & \text{otherwise} \end{cases} \end{aligned}$$

Next, we increase the size of the arrays $X(v)$ from m to the number of different weights that can be obtained by a subset of A (alternatively, B). (For example, if $|A| = 3$ and $C(a_1) = 2$, $C(a_2) = 2$, and $C(a_3) = 4$, then the weights that can be obtained are $2, 4, 2 + 4 = 6, 2 + 2 + 4 = 8$, so the size of the arrays would be 4.) Let $c[r]$ be the r 'th largest such weight. Then $X(v)[r]$ is the minimum number z , such that v can be reached from s with hops of total weight (at most) $c[r]$ of the first dog and hops of total weight z of the second dog. $X(v)[r]$ is calculated as follows:

$$X(v)[r] = \min_{(u,v) \in E} \begin{cases} X(u)[r] + w_B(u,v), & w_A(u,v) = 0 \\ X(u)[r'] + w_B(u,v), & w_A(u,v) > 0 \end{cases},$$

where $c[r'] = c[r] - w_A(u,v)$. If the number of different weights that can be obtained by a subset of A (alternatively, B) is $f(A)$ (resp., $f(B)$), then the running time is $O(m^3n^3f(A))$ (resp., $O(m^3n^3f(B))$), since the only change that affects the running time is the size of the arrays $X(v)$. We thus have

Theorem 3. *The weighted chain pair simplification problem under the discrete Fréchet distance (Weighted CPS-3F) (and its corresponding minimization problem) can be solved in $O(m^3n^3 \min\{f(A), f(B)\})$ time, where $f(A)$ (resp., $f(B)$) is the number of different weights that can be obtained by a subset of A (resp., B). In particular, if only one of the chains, say B , has points with non-unit weight, then $f(A) = O(m)$, and the running time is polynomial; more precisely, it is $O(m^4n^3)$.*

Remark 2. We presented an algorithm that minimizes $\max\{|A'|, |B'|\}$ given the error parameters $\delta_1, \delta_2, \delta_3$. Another optimization version of CPS-3F is to minimize, e.g., δ_3 (while obeying the requirements specified by δ_1, δ_2 and k). It is easy to see that Algorithm 1 can be adapted to solve this version within roughly the same time bound.

5 An Efficient Implementation

The time and space complexity of Algorithm 1 (which is $O(m^3n^3 \min\{m, n\})$ and $O(m^3n^3)$, respectively) makes it impractical for our motivating biological application (as m, n could be 500~600); see Section 7. In this section, we show how to reduce the time and space bounds by a factor of mn , using dynamic programming.

We generate all configurations of the form (a_i, a_p, b_j, b_q) , where the outermost for-loop is governed by i , the next level loop by j , then p , and finally q . When a new configuration $v = (a_i, a_p, b_j, b_q)$ is generated, we first check whether it is *possible*. If it is not possible, we set $X(v)[r] = \infty$, for $1 \leq r \leq m$, and if it is, we compute $X(v)[r]$, for $1 \leq r \leq m$.

We also maintain for each pair of indices i and j , three tables $C_{i,j}$, $R_{i,j}$, $T_{i,j}$ that assist us in the computation of the values $X(v)[r]$:

$$\begin{aligned} C_{i,j}[p, q, r] &= \min_{1 \leq p' \leq p} X(a_i, a_{p'}, b_j, b_q)[r] \\ R_{i,j}[p, q, r] &= \min_{1 \leq q' \leq q} X(a_i, a_p, b_j, b_{q'})[r] \\ T_{i,j}[p, q, r] &= \min_{\substack{1 \leq p' \leq p \\ 1 \leq q' \leq q}} X(a_i, a_{p'}, b_j, b_{q'})[r] \end{aligned}$$

Notice that the value of cell $[p, q, r]$ is determined by the value of one or two previously-determined cells and $X(a_i, a_p, b_j, b_q)[r]$ as follows:

$$\begin{aligned} C_{i,j}[p, q, r] &= \min\{C_{i,j}[p-1, q, r], X(a_i, a_p, b_j, b_q)[r]\} \\ R_{i,j}[p, q, r] &= \min\{R_{i,j}[p, q-1, r], X(a_i, a_p, b_j, b_q)[r]\} \\ T_{i,j}[p, q, r] &= \min\{T_{i,j}[p-1, q, r], T_{i,j}[p, q-1, r], X(a_i, a_p, b_j, b_q)[r]\} \end{aligned}$$

Observe that in any configuration that can immediately precede the current configuration (a_i, a_p, b_j, b_q) , the man is either at a_{i-1} or at a_i and the woman is either at b_{j-1} or at b_j (and the dogs are at $a_{p'}$, $p' \leq p$, and $b_{q'}$, $q' \leq q$, respectively). The “saving” is achieved, since

now we only need to access a constant number of table entries in order to compute the value $X(a_i, a_p, b_j, b_q)[r]$.

One can illustrate the algorithm using the matrix in Figure 3. There are mn large cells, each of them containing a matrix of size mn . The large cells correspond to the positions of the man and the woman. The inner matrices correspond to the positions of the two dogs (for given positions of the man and woman). Consider an optimal “walk” of the gang that ends at cell (a_i, a_p, b_j, b_q) (marked by a full circle), such that the first dog has visited r points. The previous cell in this “walk” must be in one of the 4 large cells $(a_i, b_j), (a_{i-1}, b_j), (a_i, b_{j-1}), (a_{i-1}, b_{j-1})$. Assume, for example, that it is in (a_{i-1}, b_j) . Then, if it is in the blue area, then $X(a_i, a_p, b_j, b_q)[r] = C_{i-1,j}[p-1, q, r-1]$ (marked by an empty square), since only the position of the first dog has changed when the gang moved to (a_i, a_p, b_j, b_q) . If it is in the purple area, then $X(a_i, a_p, b_j, b_q)[r] = R_{i-1,j}[p, q-1, r] + 1$ (marked by a x), since only the position of the second dog has changed. If it is in the orange area, then $X(a_i, a_p, b_j, b_q)[r] = T_{i-1,j}[p-1, q-1, r-1] + 1$ (marked by an empty circle), since the positions of both dogs have changed. Finally, if it is the cell marked by the full square, then simply $X(a_i, a_p, b_j, b_q)[r] = X(a_{i-1}, a_p, b_j, b_q)[r]$, since both dogs have not moved. The other three cases, in which the previous cell is in one of the 3 large cells $(a_i, b_j), (a_i, b_{j-1}), (a_{i-1}, b_{j-1})$, are handled similarly.

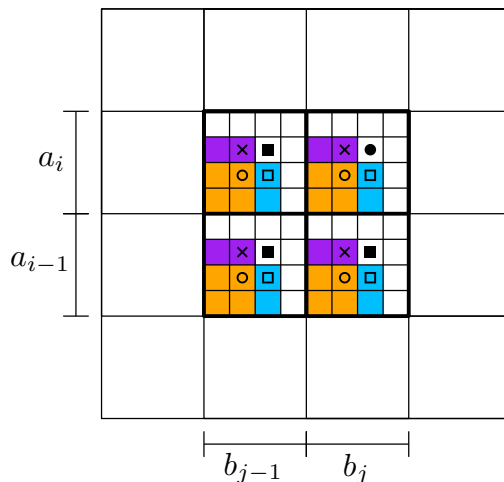


Figure 3: Illustration of Algorithm 2.

We are ready to present the dynamic programming algorithm. The initial configurations correspond to cells in the large cell (a_1, b_1) . For each initial configuration (a_1, a_p, b_1, b_q) , we set $X(a_1, a_p, b_1, b_q)[1] = 1$.

Algorithm 2 CPS-3F using dynamic programming

for $i = 1$ to m

for $j = 1$ to n

for $p = 1$ to m

for $q = 1$ to n

for $r = 1$ to m

$$X_{(-1,0)} = \min \begin{cases} C_{i-1,j}[p-1, q, r-1] \\ R_{i-1,j}[p, q-1, r] + 1 \\ T_{i-1,j}[p-1, q-1, r-1] + 1 \\ X(a_{i-1}, a_p, b_j, b_q)[r] \end{cases}$$

$$X_{(0,-1)} = \min \begin{cases} C_{i,j-1}[p-1, q, r-1] \\ R_{i,j-1}[p, q-1, r] + 1 \\ T_{i,j-1}[p-1, q-1, r-1] + 1 \\ X(a_i, a_p, b_{j-1}, b_q)[r] \end{cases}$$

$$X_{(-1,-1)} = \min \begin{cases} C_{i-1,j-1}[p-1, q, r-1] \\ R_{i-1,j-1}[p, q-1, r] + 1 \\ T_{i-1,j-1}[p-1, q-1, r-1] + 1 \\ X(a_{i-1}, a_p, b_{j-1}, b_q)[r] \end{cases}$$

$$X_{(0,0)} = \min \begin{cases} C_{i,j}[p-1, q, r-1] \\ R_{i,j}[p, q-1, r] + 1 \\ T_{i,j}[p-1, q-1, r-1] + 1 \end{cases}$$

$$X(a_i, a_p, b_j, b_q)[r] = \min\{X_{(-1,0)}, X_{(0,-1)}, X_{(-1,-1)}, X_{(0,0)}\}$$

$$C_{i,j}[p, q, r] = \min\{C_{i,j}[p-1, q, r], X(a_i, a_p, b_j, b_q)[r]\}$$

$$R_{i,j}[p, q, r] = \min\{R_{i,j}[p, q-1, r], X(a_i, a_p, b_j, b_q)[r]\}$$

$$T_{i,j}[p, q, r] = \min\{T_{i,j}[p-1, q, r], T_{i,j}[p, q-1, r], X(a_i, a_p, b_j, b_q)[r]\}$$

return $\min_{r,p,q} \max\{r, X(a_m, a_p, b_n, b_q)[r]\}$

Theorem 4. *The minimization version of the chain pair simplification problem under the discrete Fréchet distance (CPS-3F) can be solved in $O(m^2n^2 \min\{m, n\})$ time.*

6 1-Sided Chain Pair Simplification

Sometimes, one of the two input chains, say B , is much shorter than the other, possibly because it has already been simplified. In these cases, we only want to simplify A , in a way that maintains the resemblance between the two input chains. We thus define the 1-sided chain pair simplification problem.

Problem 3 (1-Sided Chain Pair Simplification).

Instance: Given a pair of polygonal chains A and B of lengths m and n , respectively, an integer k , and two real numbers $\delta_1, \delta_3 > 0$.

Problem: Does there exist a chain A' of at most k vertices, such that the vertices of A' are from A , $d_{dF}(A, A') \leq \delta_1$, and $d_{dF}(A', B) \leq \delta_3$?

The optimization version of this problem can be solved using similar ideas to those used in the solution of the 2-sided problem. Here a *possible* configuration is a 3-tuple (a_i, a_p, b_j) , where $d(a_i, a_p) \leq \delta_1$ and $d(a_p, b_j) \leq \delta_3$. We construct a graph and find a shortest path from one of the starting configurations to one of the final configurations; see Algorithm 3. Arguing as for Algorithm 1, we get that $|V| = O(m^2n)$ and $|E| = O(|V|m) = O(m^3n)$. Moreover, it is easy to see that the running time of Algorithm 3 is $O(m^3n)$, since it does not maintain an array for each vertex.

Algorithm 3 1-sided CPS-3F

1. Create a directed graph $G = (V, E)$ with a weight function w , such that:

- $V = \{(a_i, a_p, b_j) \mid d(a_i, a_p) \leq \delta_1 \text{ and } d(a_p, b_j) \leq \delta_3\}$.
- $E = \{((a_i, a_p, b_j), (a_{i'}, a_{p'}, b_{j'})) \mid i \leq i' \leq i+1, p \leq p', j \leq j' \leq j+1\}$.
- For each $((a_i, a_p, b_j), (a_{i'}, a_{p'}, b_{j'})) \in E$, set

$$w((a_i, a_p, b_j), (a_{i'}, a_{p'}, b_{j'})) = \begin{cases} 1, & p < p' \\ 0, & \text{otherwise} \end{cases}$$

- Let S be the set of starting configurations and let T be the set of final configurations.
2. Sort V topologically.
 3. Set $X(s) = 0$, for each $s \in S$.
 4. For each $v \in V \setminus S$ (advancing from left to right in the sorted sequence) do:

$$X(v) = \min_{(u,v) \in E} \{X(u) + w(u, v)\}.$$

5. Return $k^* = \min_{t \in T} X(t)$.
-

To reduce the running time we use dynamic programming as in Section 5. We generate all configurations of the form (a_i, a_p, b_j) . When a new configuration $v = (a_i, a_p, b_j)$ is generated, we first check whether it is *possible*. If it is not possible, we set $X(v) = \infty$, and if it is, we compute $X(v)$.

We also maintain for each pair of indices i and j , a table $A_{i,j}$ that assists us in the computation of the value $X(v)$:

$$A_{i,j}[p] = \min_{1 \leq p' \leq p} X(a_i, a_{p'}, b_j).$$

Notice that $A_{i,j}[p]$ is the minimum of $A_{i,j}[p-1]$ and $X(a_i, a_p, b_j)$.

We observe once again that in any configuration that can immediately precede the current configuration (a_i, a_p, b_j) , the man is either at a_{i-1} or at a_i and the woman is either at b_{j-1} or at b_j (and the dog is at $a_{p'}$, $p' \leq p$). The “saving” is achieved, since now we only need to access a constant number of table entries in order to compute the value $X(a_i, a_p, b_j)$. We obtain the following dynamic programming algorithm whose running time is $O(m^2n)$.

Algorithm 4 1-sided CPS-3F using dynamic programming

for $i = 1$ to m

for $j = 1$ to n

for $p = 1$ to m

$$X_{(-1,0)} = \min \begin{cases} A_{i-1,j}[p-1] + 1 \\ X(a_{i-1}, a_p, b_j) \end{cases}$$

$$X_{(0,-1)} = \min \begin{cases} A_{i,j-1}[p-1] + 1 \\ X(a_i, a_p, b_{j-1}) \end{cases}$$

$$X_{(-1,-1)} = \min \begin{cases} A_{i-1,j-1}[p-1] + 1 \\ X(a_{i-1}, a_p, b_{j-1}) \end{cases}$$

$$X_{(0,0)} = A_{i,j}[p-1] + 1$$

$$X(a_i, a_p, b_j) = \min\{X_{(-1,0)}, X_{(0,-1)}, X_{(-1,-1)}, X_{(0,0)}\}$$

$$A_{i,j}[p] = \min\{A_{i,j}[p-1], X(a_i, a_p, b_j)\}$$

return $\min_p \{X(a_m, a_p, b_n)\}$

Theorem 5. *The 1-sided chain pair simplification problem under the discrete Fréchet distance can be solved in $O(m^2n)$ time.*

We now study the natural problem that is obtained from the 1-sided chain pair simplification problem by omitting the requirement that $d_{dF}(A, A') \leq \delta_1$.

Problem 4 (Relaxed 1-Sided Chain Pair Simplification).

Instance: Given a pair of polygonal chains A and B of lengths m and n , respectively, an integer k , and a real number $\delta > 0$.

Problem: Does there exist a chain A' of at most k vertices, such that the vertices of A' are from A and $d_{dF}(A', B) \leq \delta$?

This problem induces two optimization problems (as in [BJW⁺08]), depending on whether we wish to optimize the length of A' or the distance between A' and B . Below we solve both of them, beginning with the former problem.

6.1 Minimizing k given δ

In this problem, we wish to minimize the length of A' without exceeding the allowed error bound.

Problem 5. Given two chains $A = (a_1, \dots, a_m)$ and $B = (b_1, \dots, b_n)$ and an error bound $\delta > 0$, find a simplification A' of A of minimum length, such that the vertices of A' are from A and $d_{dF}(A', B) \leq \delta$.

For $B = A$, Bereg et al. [BJW⁺08] presented an $O(n^2)$ -time dynamic programming algorithm. (For the case where the vertices of A' are not necessarily from A , they presented an $O(n \log n)$ -time greedy algorithm.)

Theorem 6. *Problem 5 can be solved in $O(mn)$ time and space.*

Proof. We present an $O(mn)$ -time dynamic programming algorithm. The algorithm finds the length of an optimal simplification; the actual simplification is constructed by backtracking the algorithm's actions.

Define two $m \times n$ tables, O and X . The cell $O[i, j]$ will store the length of a minimum-length simplification A^i of $A[i \dots m]$ that begins at a_i and such that $d_{dF}(A^i, B[j \dots n]) \leq \delta$. The algorithm will return the value $\min_{1 \leq i \leq m} O[i, 1]$.

We use the table X to assist us in the computation of O . More precisely, we define:

$$X[i, j] = \min_{i' \geq i} O[i', j].$$

Notice that $X[i, j]$ is simply the minimum of $X[i + 1, j]$ and $O[i, j]$.

We compute $O[-, -]$ and $X[-, -]$ simultaneously, where the outer for-loop is governed by (decreasing) i and the inner for-loop by (decreasing) j . First, notice that if $d(a_i, b_j) > \delta$, then there is no simplification fulfilling the required conditions, so we set $O[i, j] = \infty$. Second, the entries (in both tables) where $i = m$ or $j = n$ can be handled easily. In general, if $d(a_i, b_j) \leq \delta$, we set

$$O[i, j] = \min\{O[i, j + 1], X[i + 1, j + 1] + 1\}.$$

We now justify this setting. Let A^i be a minimum-length simplification of $A[i \dots n]$ that begins at a_i and such that $d_{dF}(A^i, B[j \dots n]) \leq \delta$. The initial configuration of the joint walk along A^i and $B[j \dots n]$ is (a_i, b_j) . The next configuration is either (a_i, b_{j+1}) , $(a_{i'}, b_j)$ for some $i' \geq i+1$, or $(a_{i'}, b_{j+1})$ for some $i' \geq i+1$. However, clearly $X[i+1, j+1] \leq X[i+1, j]$, so we may disregard the middle option. \square

6.2 Minimizing δ given k

In this problem, we wish to minimize the discrete Fréchet distance between A' and B , without exceeding the allowed length.

Problem 6. Given two chains $A = (a_1, \dots, a_m)$ and $B = (b_1, \dots, b_n)$ and a positive integer k , find a simplification A' of A of length at most k , such that the vertices of A' are from A and $d_F(A', B)$ is minimized.

For $B = A$, Bereg et al. [BJW⁺08] presented an $O(n^3)$ -time dynamic programming algorithm. (For the case where the vertices of A' are not necessarily from A , they presented an $O(kn \log n \log(n/k))$ -time greedy algorithm.) We give an $O(mn \log(mn))$ -time algorithm for our problem, which yields an $O(n^2 \log n)$ -time algorithm for $B = A$, thus significantly improving the result of Bereg et al.

Theorem 7. *Problem 6 can be solved in $O(mn \log(mn))$ time and $O(mn)$ space.*

Proof. Set $D = \{d(a, b) | a \in A, b \in B\}$. Then, clearly, $d_F(A', B) \in D$, for any simplification A' of A . Thus, we can perform a binary search over D for an optimal simplification of length at most k . Given $\delta \in D$, we apply the algorithm for Problem 5 to find (in $O(mn)$ time) a simplification A' of A of minimum length such that $d_F(A', B) \leq \delta$. Now, if $|A'| > k$, then we proceed to try a larger bound, and if $|A'| \leq k$, then we proceed to try a smaller bound. After $O(\log(mn))$ iterations we reach the optimal bound. \square

7 Some Empirical Results

In this section, we show some empirical results obtained by running a C++ implementation of Algorithm 2 on a standard desktop machine. The best available algorithm prior to this work was Algorithm FIND-CPS3F⁺, i.e., the algorithm (mentioned in the introduction) for the optimization version of CPS-3F⁺, proposed by Wylie and Zhu [WZ13]. This algorithm is a 2-approximation algorithm for the optimization version of CPS-3F [WZ13], and, obviously, it cannot outperform Algorithm 2 in the length of the simplification that it computes. The goal of this experimental study is thus twofold: (i) to verify that Algorithm 2 can cope with real datasets taken from the Protein Data Bank (PDB), and (ii) to examine the actual improvement obtained in the length of the simplification w.r.t. Algorithm FIND-CPS3F⁺.

Our results (summarized below) show that indeed Algorithm 2 can handle real datasets. (This is not true for the initial algorithm, i.e., Algorithm 1, whose $O(m^3n^3)$ space requirement would lead to memory overflow for most proteins. Recall that there may be as many as 500~600 α -carbon atoms along a protein backbone.) Moreover, our results show significant improvement in the length of the simplification, which is very important for the underlying structural biology applications.

As in [WZ13], we consider two cases: similar chain length and varying chain length comparisons. We use the same data and parameters as in [WZ13].

7.1 Similar chain length comparisons

We use the same seven pairs of protein backbones from the Protein Data Bank which were used in [WZ13]. To be consistent, we use the same sets of $\delta_1, \delta_2, \delta_3$ (in ångströms — note that the distance between two consecutive nodes, or α -carbon atoms, on a protein backbone, is typically between 3.7 to 3.8 ångströms). The results are summarized in Tables 1-3.

| Protein Chain(B) | $ B $ | δ_1 | δ_2 | δ_3 | $\max\{ A'' , B'' \}$ by CPS-3F ⁺ [WZ13] | $\max\{ A' , B' \}$ by CPS-3F |
|------------------|-------|------------|------------|------------|--|-----------------------------------|
| 1hfj.c | 325 | 4 | 4 | 1 | 109 | 83 |
| 1qd1.b | 325 | 4 | 4 | 21 | 126 | 82 |
| 1toh | 325 | 4 | 4 | 21 | 149 | 84 |
| 4eca.c | 325 | 4 | 4 | 6 | 111 | 83 |
| 1d9q.d | 297 | 4 | 4 | 20 | 130 | 82 |
| 4cea.b | 325 | 4 | 4 | 5 | 111 | 82 |
| 4cea.d | 325 | 4 | 4 | 5 | 113 | 84 |

Table 1: Comparison of Algorithm FIND-CPS-3F⁺ [WZ13] and Algorithm 2 in this paper with 107j.a (Chain A) of Length 325. Here A'' and B'' are the chains simplified from A and B , respectively, using the former (approximation) algorithm FIND-CPS-3F⁺, and A' and B' are the chains simplified from A and B , respectively, using Algorithm 2.

In Table 1, δ_3 is set to $\lceil d_{dF}(A, B) \rceil$. From this table one can see that with Algorithm 2, we get $\max\{A, B\} / \max\{|A'|, |B'|\} \approx 4$, while with Algorithm FIND-CPS-3F⁺, we get $\max\{A, B\} / \max\{|A''|, |B''|\} \approx 3$, using the same data and parameters. Hence, $\max\{|A'|, |B'|\} / \max\{|A''|, |B''|\} \approx 3/4$.

In Table 2, the parameters $\delta_1 = \delta_2$ are set to much larger values than in Table 1 (allowing us to set δ_3 to smaller values). The exact solutions by Algorithm 2 are even better now (w.r.t. Algorithm FIND-CPS-3F⁺). From Table 2, one can see that $\max\{|A'|, |B'|\} / \max\{|A''|, |B''|\} \approx 1/2$.

| Protein Chain(B) | $ B $ | δ_1 | δ_2 | δ_3 | $\max\{ A'' , B'' \}$ by CPS-3F ⁺ [WZ13] | $\max\{ A' , B' \}$ by CPS-3F |
|------------------|-------|------------|------------|------------|--|-----------------------------------|
| 1hfj.c | 325 | 12 | 12 | 1 | 26 | 15 |
| 1qd1.b | 325 | 15 | 15 | 12 | 21 | 11 |
| 1toh | 325 | 16 | 16 | 13 | 22 | 11 |
| 4eca.c | 325 | 12 | 12 | 3 | 27 | 16 |
| 1d9q.d | 297 | 15 | 15 | 13 | 24 | 12 |
| 4cea.b | 325 | 12 | 12 | 3 | 26 | 15 |
| 4cea.d | 325 | 12 | 12 | 3 | 32 | 16 |

Table 2: Comparison of Algorithm FIND-CPS-3F⁺ [WZ13] and Algorithm 2 in this paper with 107j.a (Chain A) of Length 325.

7.2 Varying chain length comparisons

In Table 3, we simplify A with several B chains of varying lengths. The parameter δ_1 is not set to be equal to δ_2 anymore. From the table, it can be seen that $\max\{|A'|, |B'|\} / \max\{|A''|, |B''|\}$ are mostly bounded by 2/3 to 1/2.

| Protein Chain(B) | $ B $ | δ_1 | δ_2 | δ_3 | $\max\{ A'' , B'' \}$ by CPS-3F ⁺ [WZ13] | $\max\{ A' , B' \}$ by CPS-3F |
|------------------|-------|------------|------------|------------|--|-----------------------------------|
| 3ntx.a | 322 | 10 | 10 | 5 | 39 | 25 |
| 1wls.a | 316 | 15 | 13 | 6 | 22 | 14 |
| 2eq5.a | 215 | 8 | 6 | 19 | 58 | 32 |
| 2zsk.a | 219 | 12 | 8 | 17 | 38 | 19 |
| 1zq1.a | 418 | 10 | 12 | 19 | 45 | 23 |
| 3jq0.a | 457 | 12 | 12 | 26 | 70 | 36 |
| 2fep.a | 273 | 12 | 12 | 10 | 11 | 6 |

Table 3: Comparison of Algorithm FIND-CPS-3F⁺ [WZ13] and Algorithm 2 in this paper with 107j.a (Chain A) of Length 325. Here chain B is of varying lengths.

Remark 3. Computing A', B' for a pair of protein backbones A, B might take several hours. For example, for the largest pair (i.e., $A=107j.a$ and $B=3jq0.a$ in Table 3) it takes about 20 hours, so finding heuristics for expediting the computation would be desirable. One such heuristic, is to run Algorithm FIND-CPS-3F⁺ [WZ13] to obtain a smaller upper bound on r , i.e., $\max\{|A''|, |B''|\}$ instead of m , before running Algorithm 2.

References

- [AAKS14] Pankaj K. Agarwal, Rinat Ben Avraham, Haim Kaplan, and Micha Sharir. Computing the discrete Fréchet distance in subquadratic time. *SIAM J. Comput.*, 43(2):429–449, 2014.
- [AFK⁺14] Rinat Ben Avraham, Omrit Filtser, Haim Kaplan, Matthew J. Katz, and Micha Sharir. The discrete Fréchet distance with shortcuts via approximate distance counting and selection. In *Proc. 30th Annual ACM Sympos. on Computational Geometry, SOCG’14*, page 377, 2014.
- [AG95] Helmut Alt and Michael Godau. Computing the Fréchet distance between two polygonal curves. *Internat. J. Comput. Geometry Appl.*, 5:75–91, 1995.
- [BBMM14] Kevin Buchin, Maike Buchin, Wouter Meulemans, and Wolfgang Mulzer. Four soviet walk the dog — with an application to alt’s conjecture. In *Proc. 25th Annual ACM-SIAM Sympos. on Discrete Algorithms, SODA’14*, pages 1399–1413, 2014.
- [BBW09] Kevin Buchin, Maike Buchin, and Yusu Wang. Exact algorithms for partial curve matching via the Fréchet distance. In *Proc. 20th Annual ACM-SIAM Sympos. on Discrete Algorithms, SODA’09*, pages 645–654, 2009.
- [BDS14] Maike Buchin, Anne Driemel, and Bettina Speckmann. Computing the Fréchet distance with shortcuts is NP-hard. In *Proc. 30th Annual ACM Sympos. on Computational Geometry, SOCG’14*, page 367, 2014.
- [BJW⁺08] Sergey Bereg, Minghui Jiang, Wencheng Wang, Boting Yang, and Binhai Zhu. Simplifying 3D polygonal chains under the discrete Fréchet distance. In *Proc. 8th Latin American Theoretical Informatics Sympos., LATIN’08*, pages 630–641, 2008.
- [DH13] Anne Driemel and Sarel Har-Peled. Jaywalking your dog: Computing the Fréchet distance with shortcuts. *SIAM J. Comput.*, 42(5):1830–1866, 2013.
- [EM94] Thomas Eiter and Heikki Mannila. Computing discrete Fréchet distance. Technical Report CD-TR 94/64, Information Systems Dept., Technical University of Vienna, 1994.
- [Fré06] Maurice Fréchet. Sur quelques points du calcul fonctionnel. *Rendiconti del Circolo Matematico di Palermo*, 22(1):1–72, 1906.
- [JXZ08] Minghui Jiang, Ying Xu, and Binhai Zhu. Protein structure-structure alignment with discrete Fréchet distance. *J. Bioinformatics and Computational Biology*, 6(1):51–64, 2008.

- [WLZ11] Tim Wylie, Jun Luo, and Binhai Zhu. A practical solution for aligning and simplifying pairs of protein backbones under the discrete Fréchet distance. In *Proc. Internat. Conf. Computational Science and Its Applications, ICCSA'11, Part III*, pages 74–83, 2011.
- [WZ13] Tim Wylie and Binhai Zhu. Protein chain pair simplification under the discrete Fréchet distance. *IEEE/ACM Trans. Comput. Biology Bioinform.*, 10(6):1372–1383, 2013.